

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2000-315954

(43)Date of publication of application : 14.11.2000

(51)Int.Cl.

H03M 7/30

(21)Application number : 2000-078069

(71)Applicant : LUCENT TECHNOL INC

(22)Date of filing : 21.03.2000

(72)Inventor : BENTLEY JON LOUIS  
MCILROY MALCOLM DOUGLAS

(30)Priority

Priority number : 99 273840

Priority date : 22.03.1999

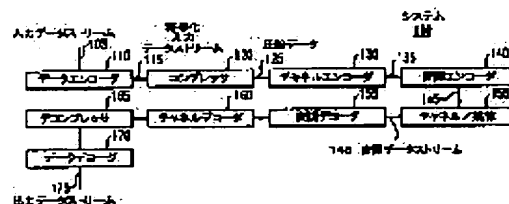
Priority country : US

(54) METHOD FOR COMPRESSING INPUTTED DATA STREAM AND DEVICE THEREFOR

(57)Abstract:

PROBLEM TO BE SOLVED: To improve data compression to a larger input data stream by crossing the input data stream and comparing an intermediate fingerprint calculated from a specified character set of the input data stream with a previously calculated and stored fingerprint.

SOLUTION: The input data stream 105 is supplied to an input data encoder 110, the input data stream is preprocessed and encoded. The encoded input data stream 115 is compressed to compressed data 125 by a compressor 120. A data stream from a channel decoder 160 is decompressed by a decompressor 165, decoded by a data decoder 170 and an output data stream 175 is created. And the input data stream is divided into groups of blocks, after that, the fingerprints are calculated for respective blocks and stored. Data compression structure to recognize correlation between different strings is further used.



## LEGAL STATUS

[Date of request for examination]

23.04.2002

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's  
decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

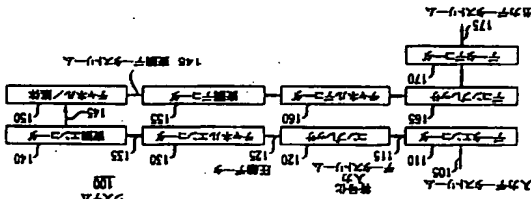
(19) 日本国特許庁 (J P) (12) 公開特許公報 (A)

(11) 特許出願公開番号  
特開2000-315954  
(P2000-315954A)  
(43) 公開日 平成12年11月14日 (2000.11.14)

(51) Int. Cl. H 03 M 7/30	識別記号	F I H 03 M 7/30	特許出願公開番号 特開2000-315954 (P2000-315954A)
(21) 出願番号 特開2000-78068 (P2000-78068)	(71) 出願人 ルセント テクノロジーズ インコーポレーテッド Lucent Technologies Inc. アメリカ合衆国 07974 ニュージャージー 一、マレービル、マウンテン アベニュー 600-700	審査請求 未請求 請求項の数33 OL (全 17 頁)	
(22) 出願日 平成12年9月21日 (2000.9.21)	(74) 代理人 井理士 三侯 弘文		
(31) 優先権主張番号 09/273840			
(32) 優先日 平成11年9月22日 (1999.9.22)			
(33) 優先権主張国 米国 (U S)			

(54) 【発明の名称】 入力データストリームの圧縮方法とその装置

(57) 【要約】  
【課題】 記憶装置の条件や伝送時間を減らすために更に大きく圧縮する方法を提供する。  
【解決手段】 圧縮プロセス (例えば、いずれかのLempel-Ziv圧縮方式) を適用する前に入力データの初期評価として入力データストリームのより長い履歴およびより長いコンテキストリングを用いる。すなわち、通常の圧縮プロセスが所望の圧縮を行うために比較的短い (最も最近の数千バイト) 入力データの履歴を用いるが、より長いコンテキストリングシーケンスを用いることにより長い履歴を用いることが組み合わさって、圧縮効率を増加させる。これは特に、繰り返し長いストリングを多く有する長い入力データストリームを圧縮する際に有効である。本発明に従うと、抽出した長コンテキストリングに対してストリングマッチングを適用することにより入力データをブリアプロセス (前処理) する。



【特許請求の範囲】

【請求項1】 入力データストリームを圧縮する方法であって、

(A) 入力データストリームを複数のデータブロックへと分割するステップと、

(B) 複数のフィンガプリントを計算するステップと、ここで、各フィンガプリントは、前記複数のデータブロックの異なる1つのデータブロックに対応し、

(C) 前記複数のフィンガプリントの各フィンガプリントは、前記複数のデータブロックの異なる1つのデータブロックに対応し、入力データストリームを前記複数のフィンガプリントと比較するステップと、

(D) 特定のフィンガプリントと入力データストリームの間にマッチングが出現すれば、そのマッチングを入力データストリーム内のそのデータストリームを圧縮するステップと、

(E) 符号化したデータストリームを圧縮データストリームへと圧縮するステップとを有することを特徴とする方法。

【請求項2】 前記入力データストリームは一連の文字からなることを特徴とする請求項1記載の方法。

【請求項3】 前記ステップ (C) は、一連の文字の個々の文字の関数として入力データストリームを順に検出し、その個々の文字のそれぞれの関数として中間フィンガプリントを計算することを特徴とする請求項2記載の方法。

【請求項4】 (F) 前記複数のフィンガプリントと共に前記複数のデータブロックをデータ構造に記憶するステップを更に有することを特徴とする請求項2記載の方法。

【請求項5】 前記ステップ (C) は、特定のフィンガプリントにマッチする一連の文字から最も長いマッチしたシーケンスを識別することを特徴とする請求項3記載の方法。

【請求項6】 前記ステップ (C) は、中間フィンガプリントを特定のフィンガプリントと比較することを特徴とする請求項3記載の方法。

【請求項7】 前記ステップ (D) は、最も長いマッチしたシーケンスの入力データストリームにおける開始位置および最も長いマッチしたシーケンスの長さの符号化することを特徴とする請求項5記載の方法。

【請求項8】 複数の文字を含む入力データストリームを圧縮コードストリームへと処理する方法であって、(A) 入力データストリームを複数のブロックに分割するステップと、ここで、各ブロックは、複数の文字のうち特定の数の文字を含み、

(B) 複数のフィンガプリントを計算するステップと、ここで、各フィンガプリントは、前記複数のデータブロックの異なる1つのデータブロックに対応し、

(C) 各フィンガプリントは複数のブロックの異なる1つのブロックに対応し、特定のフィンガプリントと特

定のブロックのいずれかの部分との間にマッチが出現するかを判断するために、複数のブロックを複数のフィンガプリントと比較するステップと、

(D) 出現したマッチそれぞれに対して入力データストリームにて識別子を符号化するステップと、ここで、この識別子は、特定のブロックのマッチ部分の入力データストリームにおける開始位置およびマッチ部分の長さを含み、

(E) 符号化された入力データストリームを圧縮コードストリームへと圧縮するステップとを有することを特徴とする方法。

【請求項9】 符号化された入力データストリームの圧縮は、Lempel-Zivコーディング技術に従って行われることを特徴とする請求項1、8記載の方法。

【請求項10】 前記ステップ (C) は、

(a) ウィンドウサイズを選択するステップと、

(b) ウィンドウサイズの関数として複数のブロックを検出するステップとを有することを特徴とする請求項8記載の方法。

【請求項11】 ウィンドウサイズは、特定の数の文字であることを特徴とする請求項10記載の方法。

【請求項12】 前記ステップ (C) は、

(c) ブロックに含まれる特定の数の文字それぞれにおけるブロックに対する中間フィンガプリントを計算するステップと、

(d) 中間フィンガプリントを特定のフィンガプリントと比較するステップとを有することを特徴とする請求項11記載の方法。

【請求項13】 前記複数のブロックの各ブロックはブロックサイズが等しいことを特徴とする請求項11記載の方法。

【請求項14】 前記ブロックサイズは10~1000文字の範囲であることを特徴とする請求項13記載の方法。

【請求項15】 前記検出 (b) は、ブロックに含まれる特定の数の文字を順に検出することを特徴とする請求項12記載の方法。

【請求項16】 (F) 前記複数のフィンガプリントと共に前記複数のデータブロックをデータ構造に記憶するステップを有し、このデータ構造は、各フィンガプリントをその異なる1つの対応するブロックと共に記憶することを特徴とする請求項8記載の方法。

【請求項17】 前記開始位置は、入力データストリームにおける文字の位置であり、マッチした部分の長さは文字の数であることを特徴とする請求項16記載の方法。

【請求項18】 (A) 入力データストリームを受信し、複数のブロックへと分割する受信器と、ここで、各ブロックは、入力データストリームからの複数の文字の特定の数の文字を含み、

(B) 複数のフィンガプリントを計算するエンコーダと、ここで、

各フィンガプリントは、前記複数のブロックの異なる1つのブロックに対応し、このエンコーダは、複数のブロックを複数のフィンガプリントと比較し、特定のフィンガプリントと特定のブロックのいずれかの部分の間でマッチが出現するかを判断し、出現した各マッチに対してそのマッチの入カデータストリームにおける識別子を符号化し、

(C) 符号化された入カデータストリームを圧縮コードストリームへと圧縮するコンプレッサとを有することを特徴とするデータ圧縮を行う装置。

【請求項19】 (D) 複数のブロックを複数のフィンガプリントと共にデータ構造に記憶するメモリを更に有し、

このデータ構造は、各フィンガプリントをその異なる1つの対応するブロックと共に記憶することを特徴とする請求項18記載のデータ圧縮を行う装置。

【請求項20】 前記複数のブロックの比較は、(a) ブロックに含まれる特定の数の各文字におけるフィンガプリントに対して中間フィンガプリントを計算し、

(b) 中間フィンガプリントを特定のフィンガプリントと比較することを特徴とする請求項18記載のデータ圧縮を行う装置。

【請求項21】 前記識別子は、特定のブロックのマッチした部分の入カデータストリームにおける開始位置およびマッチした部分の長さを含むことを特徴とする請求項19記載のデータ圧縮を行う装置。

【請求項22】 前記複数のブロックの各ブロックは、ブロックサイズが等しいことを特徴とする請求項21記載のデータ圧縮を行う装置。

【請求項23】 前記開始位置は、入カデータストリームにおける文字の位置であり、前記マッチした部分の長さは、文字の数であることを特徴とする請求項22記載のデータ圧縮を行う装置。

【請求項24】 前記コンプレッサ (C) は、圧縮コードストリームへと符号化された入カデータストリームを圧縮するのにLempel-Zivコーディングを用いることを特徴とする請求項18記載のデータ圧縮を行う装置。

【請求項25】 デジタルデータストリームを圧縮データ形式で記憶するデータ記憶システムであって、

(A) デジタルデータストリームを複数のデータブロックへと分割する手段と、ここで、各ブロックは複数の文字を有し、

(B) 複数のフィンガプリントを計算する手段と、ここで、各フィンガプリントは、前記複数のデータブロックの異なる1つのデータブロックに対応し、

(C) 複数のブロックを識別する手段と、

(D) ブロックの複数の各文字における各ブロックに対する中間フィンガプリントを計算する手段と、

(E) 中間フィンガプリントを複数のフィンガプリントの特定のフィンガプリントと比較する手段と、

(F) その特定のフィンガプリントと中間フィンガプリントの間でマッチが出現すれば、そのマッチをデジタルデータストリームにおいて符号化する手段と

(G) 符号化されたデータストリームを圧縮データ形式に圧縮する手段と、を有することを特徴とするデータ記憶システム。

【請求項26】 (H) 複数のデータブロックを複数のフィンガプリントと共にデータ構造に記憶する手段を更に有することを特徴とする請求項25記載のデータ記憶システム。

【請求項27】 前記判断する手段 (C) は、各ブロックの複数の各文字の関数としてブロックを順に判断することを特徴とする請求項26記載のデータ記憶システム。

【請求項28】 前記圧縮する手段 (G) は、符号化されたデジタルデータストリームを圧縮データ形式に圧縮するためにLempel-Zivコーディングを用いることを特徴とする請求項27記載のデータ記憶システム。

【請求項29】 前記符号化されたマッチは、特定の部分のマッチした部分のデジタルデータストリームにおける開始位置およびマッチした部分の長さを有することを特徴とする請求項25記載のデータ記憶システム。

【請求項30】 圧縮デジタル番号を処理する装置であって、

(a) 入カデジタルデータストリームを複数のデータブロックへと分割し、

(b) それぞれが複数のデータブロックの異なる1つのデータブロックに対応する複数のフィンガプリントを計算し、

(c) 入カデジタルデータストリームを複数のフィンガプリントと比較し、

(d) 特定のフィンガプリントと入カデジタルデータストリームの間でマッチが出現すれば、入カデジタルデータストリームにおいてそのマッチを符号化し、

(e) 符号化された入カデジタルデータストリームを圧縮デジタル番号に圧縮し、

(f) 圧縮デジタル番号を通信チャネルへ、供給することにより作られ、

(A) 通信チャネルから圧縮デジタル番号を受信する受信器と、

(B) 受信した圧縮デジタル番号を脱圧縮し、その脱圧縮したデジタル番号から入カデジタルデータストリームを回復するデコンプレッサとを有することを特徴とする圧縮デジタル番号を処理する装置。

【請求項31】 前記複数のブロックの各ブロックは、複数の文字を含み、前記入カデジタルデータストリームの比較 (c) は、

(g) ブロックの複数の各文字におけるブロック

に対する中間フィンガプリントを計算し、

(h) 中間フィンガプリントを特定のフィンガプリントと比較することを特徴とする請求項30記載の装置。

【請求項32】 前記複数のブロックの各ブロックは、ブロックサイズが等しいことを特徴とする請求項31記載の装置。

【請求項33】 前記複数のフィンガプリントの各フィンガプリントは対応するデータブロックの特定の文字群を要することを特徴とする請求項31記載の装置。

【請求項34】 以下の命令を含む複数の命令が記憶された機械が読み取り可能な媒体、すなわち、

(A) 機械により実行されると、入カデータストリームを圧縮コードストリームへと処理させ、

(B) 入カデータストリームを複数のブロックへと分割させ、ここで、各ブロックは複数の文字の特定の数の文字を含み、

(C) 複数のフィンガプリントを計算させ、ここで、各フィンガプリントは、前記複数のブロックの異なる1つのブロックに対応し、

(D) 複数のブロックを複数のフィンガプリントと比較させ、

(E) 特定のフィンガプリントと特定のブロックのいずれかの部分の間でマッチが出現するかどうかを判断させ、

(F) 出現した各マッチに対して、入カデータストリームにおけるそのマッチの識別子を符号化させ、ここで、この識別子は、特定のブロックのマッチした部分の入カデータストリームにおける開始位置およびマッチした部分の長さを有し、

(G) 符号化された入カデータストリームを圧縮コードストリームへと圧縮させる命令。

【請求項35】 前記複数のフィンガプリントの各フィンガプリントは、その対応するブロックの特定の文字群を要することを特徴とする請求項34記載の機械が読み取り可能な媒体。

【請求項36】 前記比較 (D) は、ブロックにおける文字の関数として行われ、ブロックに対する中間フィンガプリントは、前記複数の文字の各文字にて計算され、

前記特定のフィンガプリントと比較されることを特徴とする請求項34記載の機械が読み取り可能な媒体。

【請求項37】 前記圧縮は、Lempel-Zivコーディング技術に従って行われることを特徴とする請求項36記載の機械が読み取り可能な媒体。

【請求項38】 (A) それぞれが複数の文字からなる複数のブロックへと入カデータストリームを分割し、

(B) それぞれが前記複数のブロックの異なる1つのブロックに対応する複数のフィンガプリントを計算し、

(C) ブロックの複数の各文字の関数として入カストリームを識別し、

(D) ブロックの各文字におけるブロックに対する中間

フィンガプリントを計算し、

(E) 中間フィンガプリントを複数のフィンガプリントと比較し、

(F) 複数のフィンガプリントの特定のフィンガプリントと中間フィンガプリントの間でマッチが出現するかどうか判断し、

(G) 出現した各マッチに対して、ブロックのマッチした部分の入カデータストリームの開始位置およびマッチした部分の長さを有し、入カデータストリーム内の識別子を符号化し、

(H) 符号化された入カデータストリームを圧縮コードストリームへと圧縮し、

(I) 圧縮コードストリームを要する記憶番号を記憶媒体に供給し、

(J) 記憶番号を記憶媒体上に記憶する動作からなるプロセスを用いて作成された記憶媒体。

【発明の詳細な説明】  
[0001]

【発明の属する技術分野】 本発明は、データ記憶および通信システムに関し、特に、それらシステムの容量および利用を改善することに関する。

【従来の技術】 従来のデータ圧縮技術およびシステムはデジタルデータストリームを圧縮コードストリームへと符号化し、圧縮コードストリームを対応する元のデータストリームへとデコードして戻す。ここでコードストリームを「圧縮」としてはいるが、それは、コードストリームが通常、元のデータストリームに含まれる符号の数よりも少ない数のコードであるからである。このような小さいコードは元のデータよりも少ない量のメモリに記憶することができる。

【0003】 また、圧縮コードストリームは、圧縮されていない元のデータよりも短い時間で通信システム (例えば、有線、無線、光ファイバ通信システム) にて送信することができ、今日の通信ネットワークにおいては、コンテンツ交換の量が相対的に増大しており、データ送信および記憶容量の必要性が今までかつてないほど増えている。このように、データ圧縮は現代の通信プロトコルおよび通信ネットワークにおいて重要な役割を担っている。

【0004】 データ圧縮に有用な2つのクラスの圧縮技術は、周知のように、いわゆる特殊用途圧縮と汎用圧縮である。特殊用途圧縮技術は特殊な種類のデータを圧縮するために設計され、実装するのに比較的コストであることが多い。例えば、周知の特殊用途圧縮技術として、ランレングス符号化、ゼロパディング、ランレングス符号化、パターンの置き換えなどの技術がある。

【0005】 これら技術により、一般に、比較的小さい圧縮比となる。なぜなら、一般的に特性および冗長性を有するデータを通常圧縮するからである。圧縮比とは、

元のデータの長さに対しての圧縮コードの長さの測定値である。しかし、特殊用途圧縮技術はより一般的性質のデータ（共通な性質を多く有さないようなデータ等）を圧縮するのにあまり有効でないことが多い。

【0006】 対照的に、汎用圧縮技術は一種類のデータを特別に圧縮するのに設計されてはおらず、英数字の圧縮の際に最も有用な汎用圧縮技術は、Lempel-Zivコードである。最も有名で有用な汎用圧縮技術は、Lempel-Zivコードとマッシュ（合致）する連通したシンボルのシーケンスに基いてコードを生成することにより、データ文字のストリームを圧縮されたコードストリームへと圧縮する。各マッシュがなされるコードシンボルが生成されるに従い、このプロセスは辞書に新しいストリングエントリをも追加する。これは、データストリームにおけるマッシュしたシーケンスに加えて、データストリームにおいて遭遇する次の文字シンボルからなる。

【0011】 上述のように、Lempel-Zivコーディングの要点は、元のデータストリーム（例えば、送信されるドキュメント）において繰り返されるストリングやサブストリングを見つけて、圧縮される「テキストエントリ」における繰り返されたフレーズは元のデータストリームにおいて前に出現した場所へのポインタによって置き換えられる。従って、この方法で圧縮されるデータ（例えば、テキスト）をデコードすることには、ポインタが指すようにデコードしたテキストでポインタを置き換えることを必要とする。

【0012】 周知のように、Lempel-Zivコーディングを用いる際に設計上主に考慮することとして、ポインタがどれくらい戻るか制限を設け、その制限をどのようにするか決めることがある。別の設計上は考慮することとして所望の制限内のどのサブストリングがポインタのターゲットとするかということがある。すなわち、前のテキストへのポインタのリーチの制限がなくなったり、いわゆるグローバルウィンドウ（growing window）、あるいは前の「N」文字の固定サイズのウィンドウに制限される。ここで、Nは通常、数千文字（例えば、3キロバイト）の範囲である。

【0013】 このコーディングに従うと、ストリングの繰り返しは、両方のストリングがウィンドウに出現した場合にのみ発見され圧縮される。このようなLempel-Zivコーディングの設計時には、速さ、メモリー条件、圧縮比の面で妥協して決められる。ウィンドウをスライドすることには少なくとも1つの欠点がある。ウィンドウをスライドさせる方法は、入力テキストにおいて速くに出るストリング（例えば、10000文字分）を見つ

けることはできない。

【0014】

【発明が解決しようとする課題】 上述のようなLZ1、LZ2、LZW圧縮方式のような従来の圧縮方式は有効なデータ圧縮を提供しているが、配置装置の条件や伝送

時間を減らすために更に大きく圧縮する方式が望まれている。

【0015】

【課題を解決するための手段】 本発明は、比較的低い圧縮比を実現する方法および装置を提供する。これは、圧縮プロセス（例えば、いずれかのLempel-Ziv圧縮方式）を適用する前に入力データの初期評価として入力データストリームのより長い履歴およびより長いコンテキストストリングを用いることを我々が認識したことに基づいてい

る。すなわち、通常の圧縮プロセスが所望の圧縮を行うために比較的短い（最も最近の数キロバイト）入力データの履歴を用いるが、より長いコンテキストストリングエントリを用いることにより長い履歴を用いることが組み合わさって、圧縮効率を増加できることを我々は認識した。これは特に、繰り返す長いストリングを多く有する長い入力ストリームを圧縮する際に有効である。

【0016】 本発明に従うと、抽出した長コンテキストストリングに対してストリングマッチングを適用することになり入力データをサブプロセス（前処理）する。好ましい実施例に従うと、入力データは、個々のブロックが均等なサイズ（例えば、文字長）を有するように一連のブロックへと分割される。また、好ましい実施例に従うと、各ブロックに対していわゆるフィンガプリント（指紋）が計算され記憶される。フィンガプリントは、大きいテキストストリングの比較的小さいシグネチャである。例えば、千文字のストリングが32ビット長のフィンガプリントへとマッピングされる。従って、同一のストリングは常に同じフィンガプリントを有する。また、等しくないストリングはほとんど常に等しくないフィンガプリントを有する（特定の確率ファクタ内）。

【0017】 本発明に従うと、入力データストリームが横断され（traverse）、入力ストリームの特定の文字セットから計算した（文字ベースで）中間フィンガプリントと、前に計算し記憶したフィンガプリントとの間で比較がなされる。好ましい実施例に従うと、入力ストリームは、均等なブロックサイズを有するスライディングウィンドウの隣接として横断され、中間フィンガプリントは現在の文字ウィンドウから計算され、前に記憶したフィンガプリントと比較される。

【0018】 好ましい実施例に従い、フィンガプリントの間にマッシュを抽出すると、入力ストリームは、抽出したマッシュの関数として決められた識別子とともに符号化される。符号化された識別子は、元の入力ストリームにおけるマッシュストリングの開始位置およびストリング長を含む。その後、好ましい実施例に従うと、前述の符号化された入力ストリームに対し更にLempel-Ziv圧縮を用いて圧縮がなされる。

【0019】 本発明に従うと、全体の記憶容量の条件を余り大きくせずに、長いコンテキストストリングを識別し、入力データの大きな履歴を参照することができる。本発明に

従って、様々な圧縮方法によって大きな圧縮比を実現することができ、すなわち、本発明の原理はいずれの特定の圧縮方式には依存せず、広い範囲の圧縮方式において本発明の様々な原理を用いる利点を発揮することができ

る。

【0020】 ストリングマッチングメカニズムとしてフィンガプリントを用いることは新しくはない。テキスト処理システムにおけるストリングマッチングに対してフィンガプリントは用いられている。具体的には、テキストファイルにおいて長いコンテキストストリングを検索する際に用いられている。例えば、文献、R. M. Harper and M. O. Rabin, "Efficient Randomized Pattern-Matching Algorithms", IBM Res. Develop., Vol. 31, No. 2, pp. 249-260, March 1987は、ストリング検索の際にフィンガプリントを用いることを記載している。しかし、我々はフィンガプリントを確率な圧縮ツールとして導入できることを認識した。これにより、全体の記憶容量条件を余り大きくせずに、繰り返す長いストリングを多数有する、大きい入力データストリームに対してデータ圧縮を改善することができ

【0021】

【発明の実施の形態】 圧縮プロセス（例えば、Lempel-Ziv圧縮方式のいずれか）を適用する前に入力データの初期評価として入力データストリームのより長い履歴およびより長いコンテキストストリングを用いることの認識に基づいて、本発明は比較的圧縮比を実現する方法および装置を提供する。通常の圧縮プロセスが所望の圧縮をするために入力データの比較的短い履歴（最も最近のストリング）を用いるのが望ましくはより長いコンテキストストリングと共により長い履歴を用いることにより、特に、繰り返す長いストリングを多数有するより長い入力ストリーム（例えば、大規模データベース）を圧縮する際に、圧縮効率を増加することができることを認識した。

【0022】 本発明は、これらの方法を実現する方法および装置の形態で実施することができ、また本発明は、FD (floppy (登録商標) diskette)、CD-ROM、ハードディスクドライブ、機械が読み取り可能な記憶媒体のような具体的な媒体に実装されるプログラムコードの形態で実現することもできる。この場合、プログラムコードが機械（例えば、コンピュータ）へとロードされ機械によって実行されると、その機械が本発明を実行する装置となる。また、本発明はプログラムコードの形態で実装することができ、例えば、機械へとロードされおよび/または機械によって電磁放射のような何らかの伝送媒体によって送信されるようなプログラムコードの形態で実装される。プログラムコードが利用プロセスに実装された場合に、プログラムコードセグメントがプロセスと組合わさって、特定のロジック回路と類似するように動作するユニークなデバ

イスを与える。

【0023】図1は、本発明に従ってデータを圧縮した状態で圧縮したデータを送信するためのシステム100のブロック図である。システム100は、ほんの少しだけ名前を挙げただけである。システム100は、有線、無線、光ファイバ等によって伝送媒体（例えば、有線、無線、光ファイバ等）上に情報を送信するのにも用いられる。また、システム100は、例えば、コンピュータのディスクドライブのより一般的な磁気媒体、CD-ROMのような光学的に読み取り可能な媒体、インターネット上の媒体へと情報を記録し、またはそれらから情報を読み取るのにも有用である。

【0024】従って、本発明に従って圧縮されたデータ（例えば、圧縮されたオーディオデータ）を磁気ディスクドライブのような磁気媒体、CD-ROM等の記憶媒体に記録することができ、この記憶媒体からデータを読み出すことができる。図1において、入力データストリーム105（例えば、オーディオデータ）は、テストが入力データストリーム110に供給される。下で詳細に述べるように、本発明に従う入力データストリーム110は、抽出する良いコンストラクティングに対してストリングマッチング技術を用いることにより、本発明に従って入力データストリームをプリプロセスおの多くの符号化される。この符号化プロセスに関連する本発明の多くの原理は下で図2に示した例を特に参照してよりよく説明する。

【0025】図1に戻ると、本発明に従って作られた符号化入力データストリーム115はコンプレッサ120に渡される。好ましい実施例に従うコンプレッサ120は、符号化入力データストリーム115を圧縮データ222に圧縮する。ここで、上述のように、277コーディングを適用する。本発明に従って符号化入力データストリーム115を圧縮するのに用いられるLempel-Zivタイプの圧縮をも有効に用いて本発明の原理の利点を実現することができると考慮されたい。

【0002】次に圧縮データ125はチャネルエンコーディングされた圧縮データ120により符号化された情報121と情報を加え、エラー検出やデータ読み取りプロセスにおける訂正可能にする。伝統的にチャネル符号化技術としてでは、シンボルが1もしくは複数のデータビットで表されるシンボルのシーケンスを符号化する周知のReed-Solomon符号化がある。次にこれらシンボルは変調工法によってシンボルのシーケンスを符号化する周知の変調方式により変調符号化され、変調されたデータを送信されるかあるいは媒体150上に記録されるチャネルシーケンスを定める。

【0027】ノイズや干渉は多くの場合、データストリームの伝送や記録時にチャネル/媒体150にて投入される。従って、変調デコーダ155、チャネルデコーダ160はノイズとともに変調データストリーム145を受け取り、周知な方法で、チャネルエンコーダ130、変

図 1 の符号化プロセスをそれぞれ逆にたどると、チャネルデコーダ 16 からデータストリーム 12 が生成される。図 7、8 に開示されるように本発明に従って、このデータストリーム 12 はデコーダ 15 により脱圧縮され、データコード 17 を作る。

【0028】本発明の多くの原理は、圧縮による相当な削減や伝送効率を実現することに關連している。図2は、本発明に従ってデータを圧縮する動作を示す流れ図であり、図1のシステムにて有用である。入力データストリーム（例えば、テキストファイル）が受取られる（210）。特定の圧縮ブロックサイズ「b」が選択される（220）。これは特定の文字の数である。好ましい実施例に従って、bは20〜1000文字の範囲で選択される。入力データストリームはサイズbのブロック群に分割される（230）。その後、本発明に従って、ブロック群の各ブロックに対しフィンガプリントが計算され格納される（240）。

【0029】好ましい実施例に従って、文獻、Karpet al., *supra* に記載された技術に従ってフィンガプリン  
が計算される。カーブはストリング検索を支援する  
のにフィンガプリントを元々使った。すなわち、長さ $n$   
のストリングが長さ $m$ のサーチパターンを含むかどうか  
である。カーブはサーチパターンの $m$ の文字をポリノ  
ミアルモジュロ(多項式法処理)した大きな素数として  
解釈した。従って、得られるフィンガプリントは、例え  
ば、32ビットワードとして格納することができ、カー  
ブの技術は入カストリングを査査し、長さ $m$ の $n-m$   
+1のサブストリングのそれぞれに対し同じフィンガプ  
リントを計算する。もしこれらフィンガプリントがマッ  
チしなければ、そのサブストリングはパターンにマッ  
チしないという結論を出す。もしそれらフィンガプ  
リントがマッパすれば、そのサブストリングはパターンにマッ  
チするかどうかの更なるチェックを行う。

【0030】カープらの技術は以下のようなフィンガブリントの幾つかの有用な特性を証明した。すなわち、

- (1) フィンガブリントを迅速に計算することができること。すなわち、フィンガブリントを  $O(m)$  の時間で初期化する事ができ、 $O(1)$  の時間である位置をスライズすることにより更新することができること。
- (2) フィンガブリントは偽のマッチを得ること、すなわち、偽しくないストリングは偽しくないフィンガブリントを持つことと対称に偽かであること。(2つの等しくないストリングが同じ3.2ビットフィンガブリントを有する確率は約  $2^{-32}$  である。)

(3) 大きな素数をランダムで選ぶことができ、テキストストリング検索においてランダム化したアルゴリズムを得ることができること。

に大きな入カデータストリームに対し増強したデータ圧縮を行うようなエレガントな圧縮ツールを導入することによってフィンガプリントを用いることを認識した。本発明に従って、全体の記憶容量の必要条件の余り増加させないで最もコモンストリングを識別して入カデータに大きな問題を課することができ、すなわち、本発明に従って、異なるストリングの間の相関を認識するデータ圧縮構成を用いる。具体的には、特定のテキストストリングの第2出現を繰り返り返しとして認識し、この第2出現を符号化した第1ストリングへの参照によって置き換える。従って、本発明に従うと、多くの圧縮方法に対し大きな圧縮比を達成することができ、

【0.0.3.2】特定の圧縮方式におけるストリングの繰り返し認識は知られている。例えば、文獻 J. G. Cleary et al., "Unbounded length context for PPM", *Computational Journal*, 40, 2/3, pp. 67-75, 1997, C. G. Nevill et al., "A space-economical suffix tree construction algorithm", *Journal of Artificial Intelligence Research*, pp. 67-82, September 1997)にはストリングの繰り返しを認識している特定の圧縮方式を記載している。しかし、本発明とは対照的にこれらの従来の方式は大きな量のメモリを必要とする。すなわち、フライングプリントを用いない方法において、n文字のファインガプリントをするのに約nのワードを主メモリに必要とする。下で議論するように、記憶条件の必要条件を激的に増加せず、本発明に従うような圧縮比を実現することができる。

【0033】図2において、好ましい実施例に従うと、各ブロックに対して算出したフィンガプリントを格納する(240)。入カストリームの各ブロック境界にてフィードバックが記憶される。また、ブライバードデータ構造がbバイトの重なり合わないブロックのそれぞれのフィンガプリントを記憶する。すなわち、好ましい実施例に従うと、各バイト1...b、b+1...2b、2b+1...3b毎々のようにフィンガプリントが記憶される。

【0034】好ましい実施例に従うと、約  $n/6$  のフィ  
ンガプリントが記憶される。この  $n$  は上述のスリーマン  
の長さである。本発明に従うと、元の入カストリーマン  
の少ない割合のみが記憶され、従って、記憶装置の必要条  
件を低減させることが可能である。また、好ましい実施例  
に従うと、入カデータストリームにおけるそのシーケン  
スの位置を符号と共にお互い互いをハッシュテーブル  
(周知なデータ構造)にてフィンガプリントを記憶し表  
現する。

【0035】図3は、図2の動作によって計算されたフ  
ィンガブリントを記憶するデータ構造300を示す。デ  
ータ構造300は入カストリームの各ブロックを記憶す  
る。例えば、ブロック1～miはブロック305～325  
として示した。それに加え、各ブロックに対し、計算し

たフィンガブリットを配座する。例えば、図3において F P 1 ~ F P n をフィンガブリット 3.30 ~ 3.50 とし示した。更に、下で述べるように、入力データストリーム（データ幅は 3.00）を横断するのにスライディングウィンドウ 3.55 を用い、上述のマッチを検出するに現在の文字のウィンドウを配座されたフィンガブリットと比較する。

【0036】より詳細には、図2において、入力データストリームを抽出し、入力文字と配列されたフィンガプリントの間にはずれを行い(250)、マッチを抽出する(260)。すなわち、入力データストリームを抽出するのにはスライディングウィンドウ(例えば、スライディングウィンドウ355)を用い、現在の文字のウィンドウに対していわゆる「中間(interlirid)」フィンガプリントを計算する。これは図において現在の文字のウィンドウで文字毎のペースで行われている。

【003.7】これら現在の文字のウィンドウにわたって計算される中間フィンガプリントはマッチを検出するため、単純フィンガプリントと比較される。具体的には、入力テキストを走査するに従って、モメンティックグリッドを見つけるためにハッシュテーブルが用いられ、コメントストリングの位置を判断する。もしマッチが出現すれば、シーケンス  $\langle \text{start}, \text{length} \rangle$  を用いてそのマッチを符号化する (2.7.0)。ここで、start は初期位置であり、length はコメントシーケンスのサイズである。【003.8】例えば、以下の入力データストリームを考  
えてみる。

The Constitution of the United States, PREAMBLE We the people of the United States, in order to form a more perfect union...

上述のような本発明の原理を適用すると、以下の符号化データストリームの結果を得ることができる。

The Constitution of the United States. PREAMBLE W  
e, the people (16, 21), in order to form a more per  
fect union...

この符号化データストリームから、繰り返しストリングは、the people of the United States は上述のように本発明の「`<start, length>`」シーケンスフォーマットを用いる識別子で符号化されていることがわかる。

【0039】好ましい実施例に従うと、マッチするフィ  
ンガプリントを有するブロックが、いわゆる偽マッチでは  
ないことを識別するため、可能な限りであるが、母が  
ーム文字よりともしくない範囲で、（入力データストリ  
ームにわたって）逆方向および順方向でマッチの拡張を  
行う。もし幾つかのブロックが現フィンガプリントとマ  
ッチすれば、そのようなブロックの中で最大のマッチ  
が本発明に従って符号化される。

【0040】フィンガプリントに対しての入力文字（中間フィンガプリント）の比較は入カストリームの終わりに到達するまで継続する（280）。その後、符号化

データストリーム（本発明に従って符号化された元の入力カデータストリーム）が周知の圧縮アルゴリズムのいずれか（例えば、Lempel-Ziv圧縮）を圧縮される（290）。

（0041）以下の導コードは、上述のような本発明の原理に従うフィンガリングの比較および符号化を記したものである。変数ipは、フィンガリングを要し、関数checkformalchはハッシュテーブルにおけるフィンガリングをルックアップし、マッチを見つけたればそのマッチを符号化する。

```
コード1
initialize fp
for (i=0; i<n; i++)
    if (!x[i] == 0)
        store (fp, i)
    a[i] を含み a[i-1] を含まないよう ip を更新
    checkformalch (fp, i)
```

（0042）上の導コードは、本発明を要するプロセスにおける実行をするために、多くのプログラム（例えば、C書庫プログラム）にて開発するのに用いることができる。例えば、図9～10は、本発明に従ってデータ圧縮するC書庫ソースプログラム900を示す。ソースプログラム900は、このソースコードプログラム900全体で用いる特定のデータ型、変数、データ構造を規定するプログラム命令を含むプログラムソースコード部分910を含む。

（0043）プログラムソースコード部分920は、上述のように本発明に従ってストリングマッチング動作を実行するプログラム命令を含む。プログラムソースコード部分930は、上述のように本発明に従って計算されるフィンガリングを保持するのに用いるハッシュテーブルデータ構造を定義するのを要するプログラム命令である。プログラムソースコード部分940、950（図10）は、上述のように本発明に従って符号化データファイルを作るために圧縮を完成させるプログラム命令を含む。

（0044）C書庫ソースプログラム900は、生来的に例示的であり本発明を理解するのを要するために要現である。本発明を具現化する他のプログラムを本発明の範囲から外れずに開発することは当業者であればできるであろう。

（0045）上述の本発明の多くの原理を更に説明するため、図4は、一連の短い非圧縮入力カデータストリーム400を、本発明に従って符号化された対応する一連の符号化出力カデータストリーム410と共に示してある。図4を調べることに、入力カデータストリーム420～450はそれぞれ、符号化出力460～490それぞれに示すようにマッチングストリングの符号化された表現によって処理されている。

（0046）また、図5、6には本発明を更に説明する

ために、大きな入力カデータファイル、具体的には、アメリカ合衆国の憲法に対する本発明の原理の適用を示している。図5は、憲法の選択部分からなる入力カデータファイル500を示している。本発明に従う憲法のテキストの圧縮は、同じ長いストリングが頻りに現れることを考へると特に本発明の利点を発揮できる。例えば、入力カデータファイル500のテキスト部分510、520、530は横つきの長い繰り返りストリングを含む。

（0047）従って、上述の本発明の多くの原理を適用すると、図6に示した符号化データファイル600を得る。600は符号化部分610、620、630を含む。これらそれぞれは入力カデータファイル500のテキスト部分510、520、530に対応する。符号化部分610、620、630はそれぞれ、エンコーディング（例えば、エンコーディング635～690）を含む。これらは、本発明に従って導かれ、ブロックサイズ＝20を用い、更に圧縮比を増しデータ伝送レートをより効率的にするために更に圧縮されることができ、（0048）例えば、エンコーディング635は、文字位置391に始まる47文字のマッチングストリング（すなわち、マッチングストリング）が抽出され符号化されたことを示す。同様に、例えば、エンコーディング675は、文字位置2439に始まる103文字のマッチングストリングが抽出され符号化されたことを示す。図9（91）圧縮プロセスを適用する前に入力データの初期評価（プリプロセス）においてフィンガリングを用いることを説明することにより、入力カデータストリームの長い問題および長いコメントストリングを用いることに基づいて比較的低い圧縮比を実現することができる。重要なことに、本発明は配電装置の必要条件を余り増加させず、いずれの特定の圧縮技術にも依存しない。すなわち、本発明は、相当な圧縮比を実現する多様な周知の圧縮アルゴリズムを用い、配電装置の必要条件や伝送時間を減らすことができる。

（0050）圧縮テスト結果を記す前に、本発明を用いて得られる符号化データファイルの圧縮を議論する。図7は、本発明の更なる原理に従ってデータを圧縮する動作700の流れ図である。本発明による符号化データファイル（例えば、符号化データファイル600）から、この図では文字毎のベースで個々の文字「c」を取り出す（710）。

（0051）特定の文字cがシンボル「<」とマッチするかどうかの判断を行う（720）。マッチしなければ、脱圧縮プロセスに従って出力カデータファイルへ文字cが書き込まれる（730）。符号化データファイルから次の文字が取り出される。マッチすれば、符号化データファイルから次の文字が取り出される（740）。再び、特定の文字cがシンボル「<」とマッチするかどうかの判断をする（750）。マッチすれば、出力カデータファイルにシンボル「<」が書き込まれる（760）。次の文字が取り

本発明の原理と組み合わせる周知の「gzip」圧縮アルゴリズム（周知のGNUによる277の実装であるgzip）を適用した結果を示す。「com|gzip %」（850）は、列820のサイズに対する列840の割合を示す。

例えば、列820のサイズを適用する元、元のファイルサイズを19.7%の割合で減らすことができた。図よりブロックサイズが減ると、最適でない点（効率的な圧縮にはブロックサイズが小さすぎ）に到達するまで圧縮の度合いは増えている。また、「total %」の見出しを有する列860は、元のファイルと比較した列840の割合を示し、aの最適な選択は、com|gzipのファイルサイズと、bの最適な選択は、com|gzipのファイルサイズの減少量を超えて22.5%のファイルサイズの減少量となっているところの31%である。

（0058）本明細書における説明は本発明の原理の例を示したのみである。当業者は、本発明の範囲を外れずに多くの構成を考へることができであろう。特に、多くの構成を考へることができであろう。特に、許容範囲の記載において特定の機能を実行する手段として表現したいずれの要素も、その機能を実行するいずれの機能をも表すように意図してある。例えば、（a）その機能を実行する回路要素の組み合わせ、あるいは（b）その機能を実行するためのソフトウェアを実行する適切な回路と組み合わせるいずれの形態におけるソフトウェア（従って、ファームウェア、オブジェクトコード、マイクロコード等を含む。）を含む。

【図面の簡単な説明】

【図1】本発明に従ってデータを圧縮し展開するシステム。

【図2】本発明に従ってデータを圧縮する動作の流れ図。これは図1のシステムにて有用である。

【図3】図2の動作に従って計算された入力カデータストリームとフィンガリングを記憶するデータ構造の例。

【図4】本発明に従って符号化された一連の符号化データストリームとともに一連の非圧縮データストリームを示す。

【図5】入力カデータファイルの選択部分。

【図6】本発明の原理に従って図5の入力カデータファイルから符号化された符号化データファイルの選択部分。

【図7】本発明に従って展開する動作の流れ図。

【図8】本発明に従ってテキストファイルを選択した圧縮結果。

【図9】図2に示すように本発明に従ってデータを圧縮するC書庫ソースコードプログラム。

【図10】図9と同様。

【図11】図7に示すように本発明に従ってデータを展開するC書庫プログラム。

【符号の説明】

100 55 本発明の多くの原理および利点を説明するため、非常に大きなファイルに関連して本発明を適用し、そのファイルの圧縮の多くの結果と比較した。サンプルとしては、周知のCD-ROMで「Project Gutenberg Compact Disc」、Walnut Creek CDROM、Walnut Creek, CAに含める全てのテキストファイルを選択させて用いた。このCD-ROMには1994の文書が含まれている。この試験のため、我々は周知のUNIX（登録商標）の「cat」コマンドを用いて全てのテキストファイルを選択した。具体的なUNIXコマンド文字列は、「cat /s/a/ixl > gull94all.ixl」であった。このように連結することにより、661220bytesの入力ファイルを得て、これに本発明の原理を適用した。

（0056）図8は、連結した当該テキストファイルを圧縮した結果の比較800を示す。ファイルサイズは元で示してあり、圧縮による元ファイルと割合をパーセントで示した。ブロックサイズb（810）を最大化させた効果を比較した。「com」の見出しを有する列820は、本発明を適用した効果を示してあり、入力カデータのファイルのサイズの変化を示してある。「com %」の見出しを有する列830は、ブロックサイズを最大化させて調べた圧縮パーセンテージを示す。

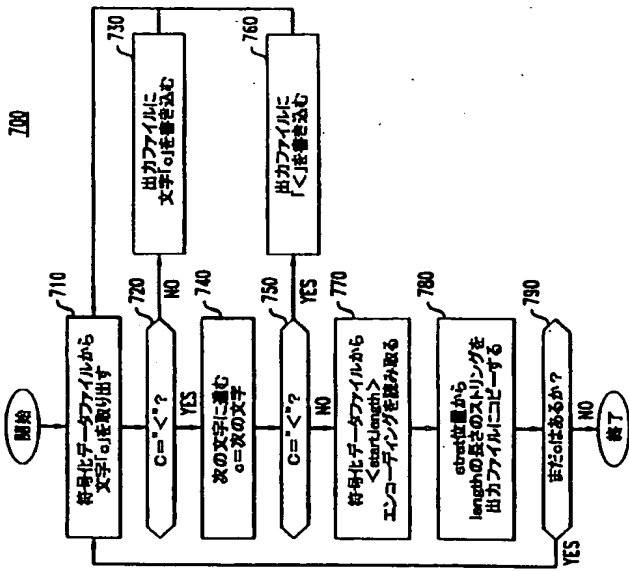
（0057）「com b |gzip」の見出しを有する列840および「com|gzip %」の見出しを有する列850は、



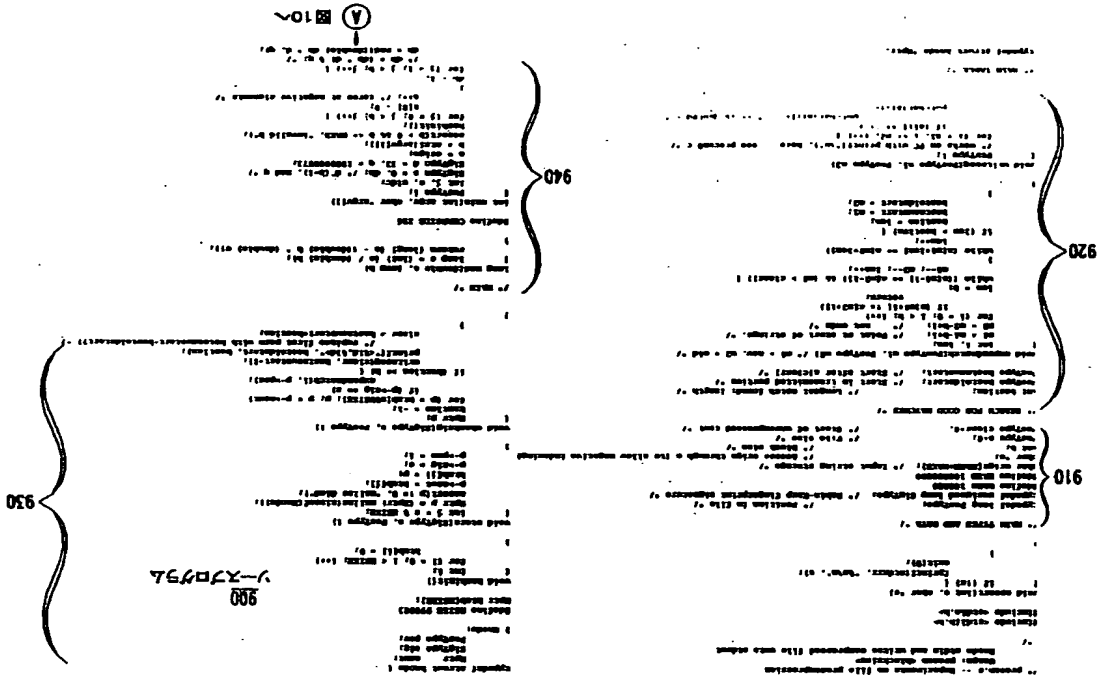




【図7】



【図9】



[図 11]

## 1100 ソースコードプログラム部分

```

/* decom.o -- decompress files made by precomp.o
   Usage: precomp
   Reads stdin and writes decompressed file onto stdout
*/

#include <stdio.h>
#include <stdlib.h>
#include <ctype.h>

typedef long PostType;
#define MAX 100000000
char a[MAX]; /* Output string storage */
PostType n=0; /* File size */

void addchar(char c)
{
    a[n++] = c;
    putchar(c);
}

int getint() /* gobble nondigit at end */
{
    int c, n = 0;
    for (;;) {
        c = getchar(); /* danger: no EOF */
        if (!isdigit(c)) break;
        n = 10*n + c - '0';
    }
    return n;
}

int main()
{
    PostType i;
    int c, start, len;
    while ((c = getchar()) != EOF) {
        if (c == '<') {
            c = getchar(); /* danger: no EOF */
            if (c == '<') { /* quoted < */
                addchar('<');
            } else {
                ungetc(c, stdin);
                start = getint();
                len = getint();
                for (i = start; i < start+len; i++)
                    addchar(a[i]);
            }
        } else
            addchar(c);
    }
    return 0;
}

```

## フロントページの続き

(71) 出願人 596077259

600 Mountain Avenue,  
Murray Hill, New Je  
rsey 07974-0636 U. S. A.

(72) 発明者 ジョン ルイス ベントリー

アメリカ合衆国、07974 ニュージャージー  
ー、ニュープロビデンス、セントラル ア  
ベニュー 449、アパートメント 2

(72) 発明者 マルコルム ダグラス マクローイ  
アメリカ合衆国、03750 ニューハンプシ  
ヤー、エトナ、ハイフィールド ロード  
1